



**ICLR**



清华大学  
Tsinghua University

# Exploring Memorization in Adversarial Training

**Yinpeng Dong<sup>1,2</sup>, Ke Xu<sup>3</sup>, Xiao Yang<sup>1</sup>, Tianyu Pang<sup>1</sup>, Zhijie Deng<sup>1</sup>, Hang Su<sup>1</sup>, Jun Zhu<sup>1,2</sup>**

<sup>1</sup> Dept. of Comp. Sci. and Tech., Institute for AI, Tsinghua University <sup>2</sup> RealAI <sup>3</sup> CMU

[{dongyinpeng, suhangss, dcszj}@mail.tsinghua.edu.cn](mailto:{dongyinpeng, suhangss, dcszj}@mail.tsinghua.edu.cn); [kx1@andrew.cmu.edu](mailto:kx1@andrew.cmu.edu);

[{yangxiao19, pty17, dzj17}@mails.tsinghua.edu.cn](mailto:{yangxiao19, pty17, dzj17}@mails.tsinghua.edu.cn)

# Adversarial Examples

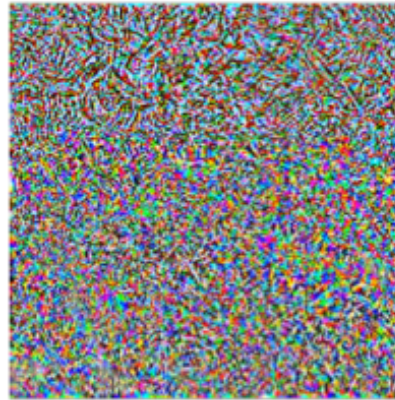


Clean images



Alps: 94.39%

Adversarial noise



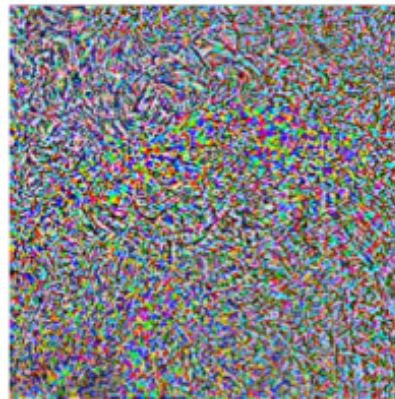
Adversarial examples



Dog: 99.99%



Puffer: 97.99%



Crab: 100.00%

(Figure is from Dong et al. 2018)

# Adversarial Training

- From the optimization view, adversarial training (AT) can be formulated as a minimax optimization problem (Madry et al., 2018)

Outer minimization: train a robust classifier

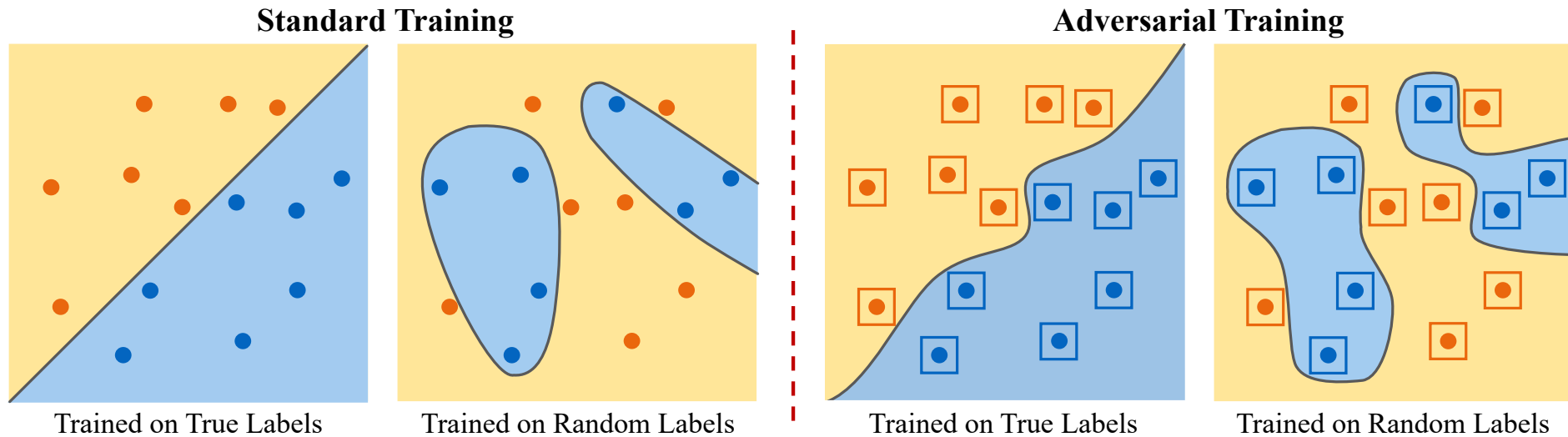
$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \max_{\delta_i \in S} L(f_{\theta}(x_i + \delta_i), y_i) \quad \leftarrow S = \{\delta: \|\delta\|_{\infty} \leq \epsilon\}$$

Inner maximization: generate an adversarial example

- Solve the inner maximization by projected gradient descent

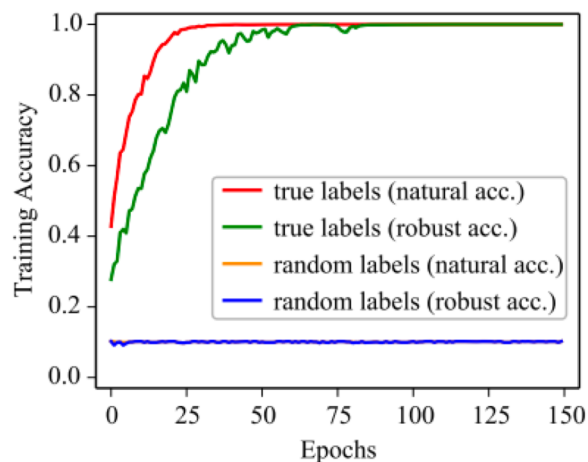
$$\delta_i^{t+1} = \Pi_S \left( \delta_i^t + \alpha \cdot \text{sign} \left( \nabla_x L(f_{\theta}(x_i + \delta_i^t), y_i) \right) \right)$$

# Memorization

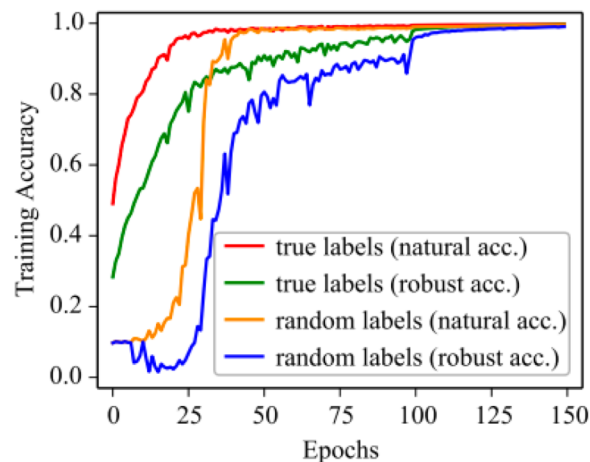


**Goal:** To facilitate a deeper understanding of model capacity, training convergence, robust generalization, and robust overfitting of the adversarially trained models.

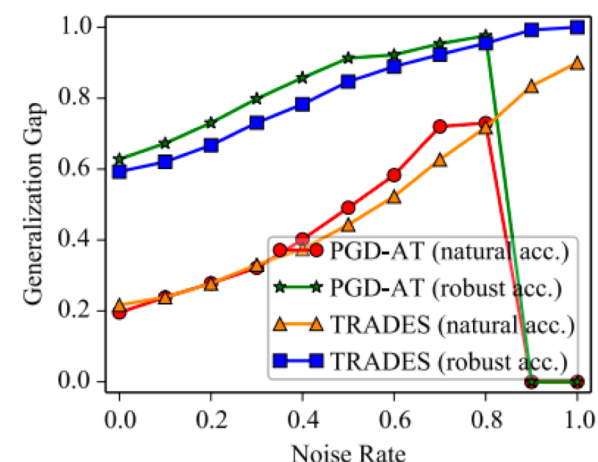
# AT with random labels



(a) PGD-AT



(b) TRADES



(c) Generalization gap growth

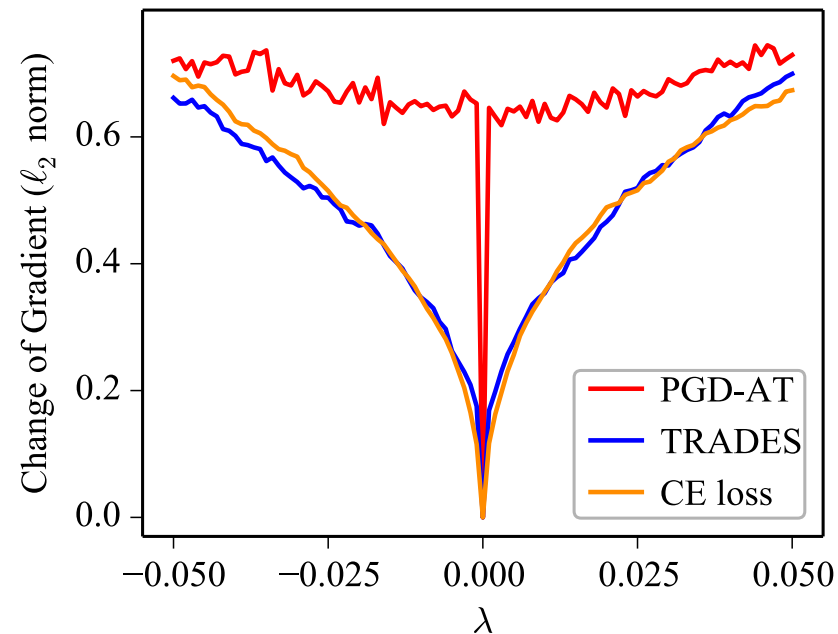
**Finding:** PGD-AT cannot converge with random labels, but TRADES can.

- This finding holds with different training settings, including network architecture, attack steps, optimizer, perturbation budget, regularizations

# Convergence Analysis

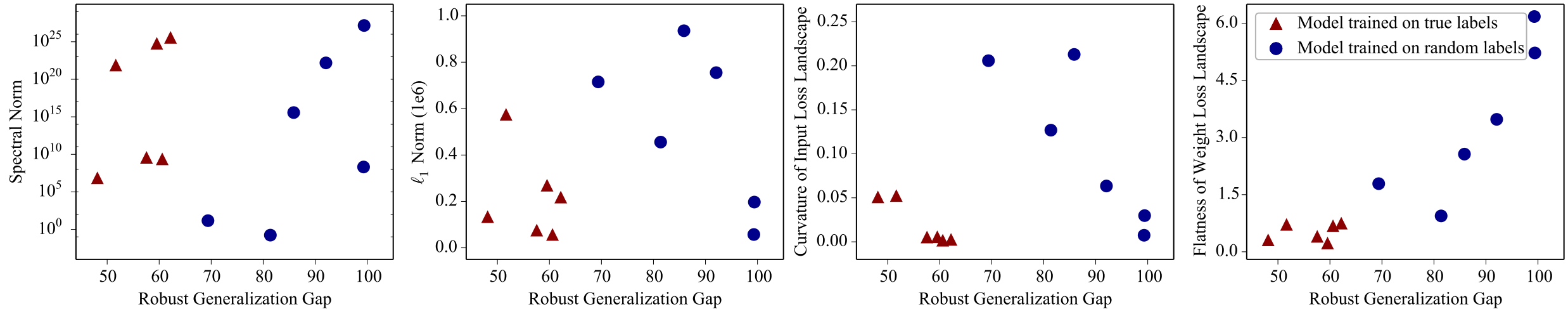


- **Gradient instability issue:** the gradient of the adversarial loss in PGD-AT changes more abruptly than TRADES.



We measure the L2 distance between the gradient at  $\theta$  and  $\theta + \lambda d$

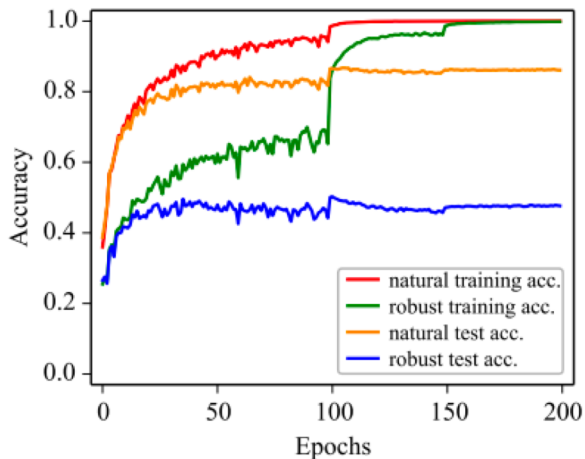
# Generalization Analysis



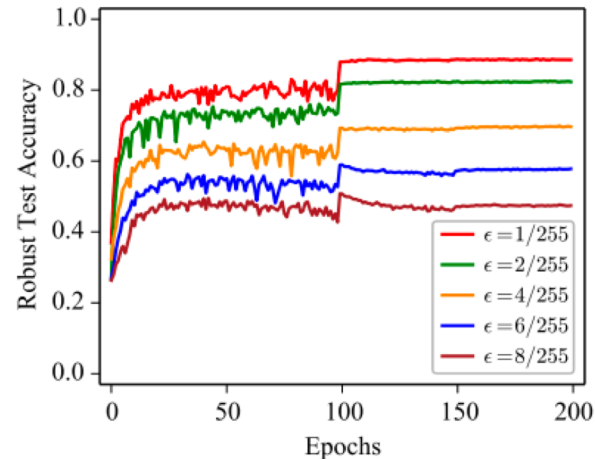
- We consider two norm-based measures and two sharpness/flatness-based measures.
- **Finding:** None of them can adequately explain and ensure robust generalization.

# Robust Overfitting Analysis

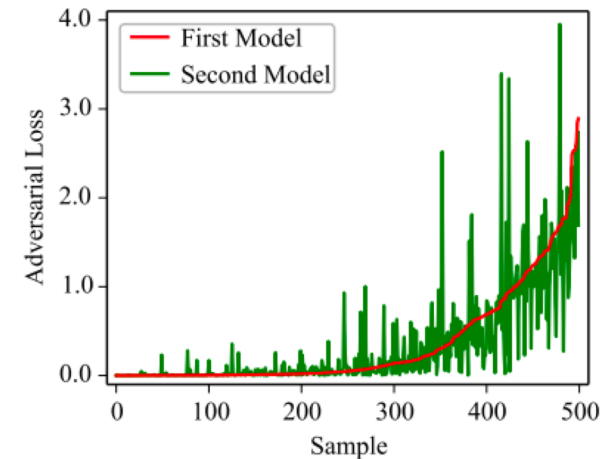
- **Argument:** robust overfitting is caused by excessive memorization of (noisy) one-hot labels.



(a)



(b)



(c)

1. Robust overfitting does not occur with a smaller perturbation budget (e.g.,  $\epsilon = 1/255$ ).
2. The hard examples are consistent across different networks.

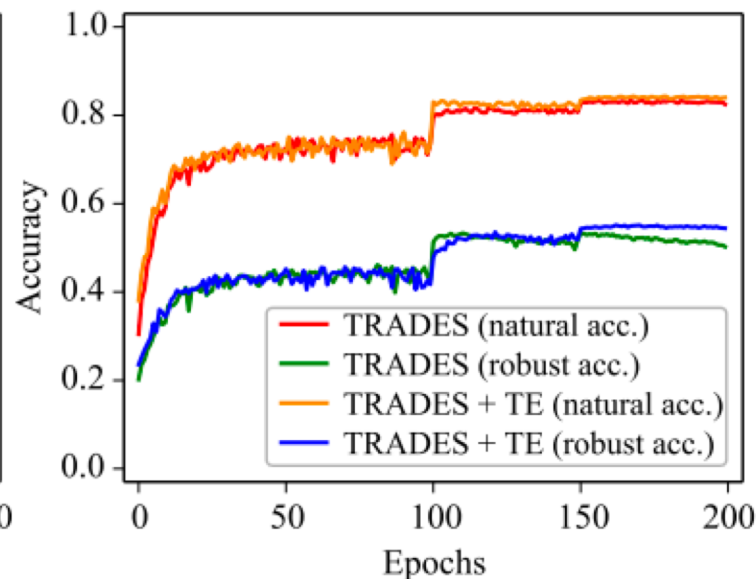
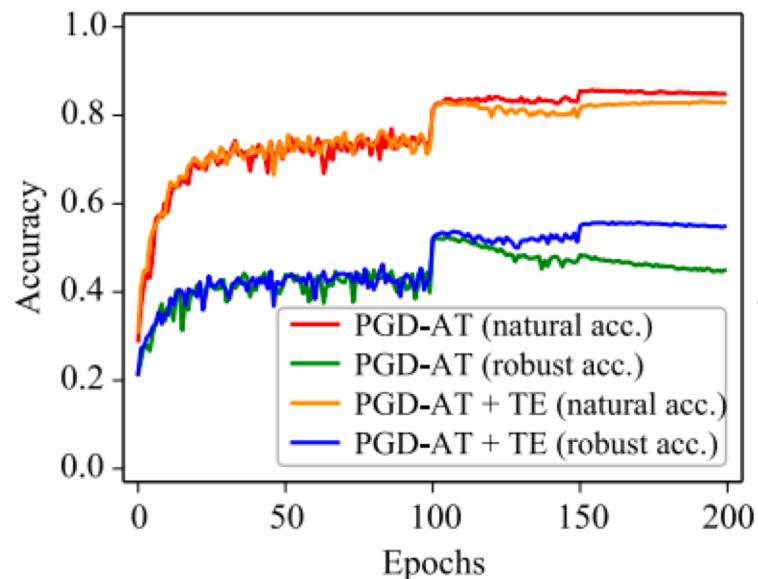


# Mitigation Algorithm

- Incorporate the **Temporal Ensembling (TE)** approach into AT

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \max_{\delta_i \in \mathcal{S}} \{L(f_{\theta}(x_i + \delta_i), y_i) + w \cdot \|f_{\theta}(x_i + \delta_i) - \hat{p}_i\|_2^2\}$$

- $\hat{p}_i = \frac{p_i}{\|p_i\|_2}, p_i = \eta p_i + (1 - \eta)f_{\theta}(x_i)$



Learning curves  
on CIFAR-10

# Empirical Results



Method	Natural Accuracy			PGD-10			PGD-1000			C&W-1000			AutoAttack		
	Best	Final	Diff	Best	Final	Diff	Best	Final	Diff	Best	Final	Diff	Best	Final	Diff
PGD-AT	<b>83.75</b>	<b>84.82</b>	-1.07	52.64	44.92	7.72	51.22	42.74	8.48	50.11	43.63	7.48	47.74	41.84	5.90
PGD-AT+TE	82.35	82.79	<b>-0.44</b>	<b>55.79</b>	<b>54.83</b>	<b>0.96</b>	<b>54.65</b>	<b>53.30</b>	<b>1.35</b>	<b>52.30</b>	<b>51.73</b>	<b>0.57</b>	<b>50.59</b>	<b>49.62</b>	<b>0.97</b>
TRADES	81.19	82.48	-1.29	53.32	50.25	3.07	52.44	48.67	3.77	49.88	48.14	1.74	49.03	46.80	2.23
TRADES+TE	<b>83.86</b>	<b>83.97</b>	<b>-0.11</b>	<b>55.15</b>	<b>54.42</b>	<b>0.73</b>	<b>53.74</b>	<b>53.03</b>	<b>0.71</b>	<b>50.77</b>	<b>50.63</b>	<b>0.14</b>	<b>49.77</b>	<b>49.20</b>	<b>0.57</b>

(a) The evaluation results on **CIFAR-10**.

Method	Natural Accuracy			PGD-10			PGD-1000			C&W-1000			AutoAttack		
	Best	Final	Diff	Best	Final	Diff	Best	Final	Diff	Best	Final	Diff	Best	Final	Diff
PGD-AT	<b>57.54</b>	<b>57.51</b>	<b>0.03</b>	29.40	21.75	7.65	28.54	20.63	7.91	27.06	21.17	5.89	24.72	19.34	5.38
PGD-AT+TE	56.45	57.12	-0.67	<b>31.74</b>	<b>30.24</b>	<b>1.50</b>	<b>31.27</b>	<b>29.80</b>	<b>1.47</b>	<b>28.27</b>	<b>27.36</b>	<b>0.91</b>	<b>26.30</b>	<b>25.34</b>	<b>0.96</b>
TRADES	57.98	56.32	1.66	29.93	27.70	2.23	29.51	26.93	2.58	25.46	24.42	1.04	24.61	23.40	1.21
TRADES+TE	<b>59.35</b>	<b>58.72</b>	<b>0.63</b>	<b>31.09</b>	<b>30.12</b>	<b>0.97</b>	<b>30.54</b>	<b>29.45</b>	<b>1.09</b>	<b>26.61</b>	<b>25.94</b>	<b>0.67</b>	<b>25.27</b>	<b>24.55</b>	<b>0.72</b>

(b) The evaluation results on **CIFAR-100**.

Method	Natural Accuracy			PGD-10			PGD-1000			C&W-1000			AutoAttack		
	Best	Final	Diff	Best	Final	Diff	Best	Final	Diff	Best	Final	Diff	Best	Final	Diff
PGD-AT	89.00	90.55	-1.55	54.51	46.97	7.54	52.22	42.85	9.37	48.66	44.13	4.53	46.61	38.24	8.37
PGD-AT+TE	<b>90.09</b>	<b>90.91</b>	<b>-0.82</b>	<b>59.74</b>	<b>59.05</b>	<b>0.69</b>	<b>57.71</b>	<b>56.46</b>	<b>1.25</b>	<b>54.55</b>	<b>53.94</b>	<b>0.61</b>	<b>51.44</b>	<b>50.61</b>	<b>0.83</b>
TRADES	<b>90.88</b>	<b>91.30</b>	<b>-0.42</b>	59.50	57.04	2.46	52.78	50.17	2.61	52.76	50.53	2.23	40.36	38.88	1.48
TRADES+TE	89.01	88.52	0.49	<b>59.81</b>	<b>58.49</b>	<b>1.32</b>	<b>58.24</b>	<b>56.66</b>	<b>1.58</b>	<b>54.00</b>	<b>53.24</b>	<b>0.76</b>	<b>51.45</b>	<b>50.16</b>	<b>1.29</b>

(c) The evaluation results on **SVHN**.



# Thanks for listening

Code is available at:

<https://github.com/dongyp13/memorization-AT>

